

Hindawi Publishing Corporation  
EURASIP Journal on Audio, Speech, and Music Processing  
Volume 2011, Article ID 426792, 17 pages  
doi:10.1155/2011/426792

## Research Article

# Phoneme and Sentence-Level Ensembles for Speech Recognition

Christos Dimitrakakis<sup>1</sup> and Samy Bengio<sup>2</sup>

<sup>1</sup>FIAS, Ruth-Moufang-Strß 1, 60438 Frankfurt, Germany

<sup>2</sup>Google, 1600 Amphitheatre Parkway, B1350-138, Mountain View, CA 94043, USA

Correspondence should be addressed to Christos Dimitrakakis, [christos.dimitrakakis@gmail.com](mailto:christos.dimitrakakis@gmail.com)

Received 17 September 2010; Accepted 20 January 2011

Academic Editor: Elmar Nöth

Copyright © 2011 C. Dimitrakakis and S. Bengio. This is an open access article distributed under the Creative Commons Attribution License, which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited.

We address the question of whether and how boosting and bagging can be used for speech recognition. In order to do this, we compare two different boosting schemes, one at the phoneme level and one at the utterance level, with a phoneme-level bagging scheme. We control for many parameters and other choices, such as the state inference scheme used. In an unbiased experiment, we clearly show that the gain of boosting methods compared to a single hidden Markov model is in all cases only marginal, while bagging significantly outperforms all other methods. We thus conclude that bagging methods, which have so far been overlooked in favour of boosting, should be examined more closely as a potentially useful ensemble learning technique for speech recognition.

## 1. Introduction

This paper examines the application of ensemble methods to hidden Markov models (HMMs) for speech recognition. We consider two methods: bagging and boosting. Both methods feature a fixed mixing distribution between the ensemble components, which simplifies the inference, though it does not completely trivialise it.

This paper follows up on and consolidates previous results [1–3] that focused on boosting. The main contributions are the following. Firstly, we use an unbiased model testing methodology to perform the experimental comparison between the various different approaches. A larger number of experiments, with additional experiments on triphones, shed some further light on previous results [2, 3]. Secondly, the results indicate that, in an *unbiased* comparison, at least for the dataset and features considered, bagging approaches enjoy a significant advantage to boosting approaches. More specifically, bagging consistently exhibited a significantly better performance than either any of the boosting approaches examined. Furthermore, we were able to obtain state-of-the-art results on this dataset using a simple bagging estimator on triphone models. This indicates that perhaps a shift towards bagging and perhaps, more generally, empirical Bayes methods may be advantageous for any further advances in speech recognition.

Section 2 introduces notation and provides some background to speech recognition using hidden Markov models. In addition, it discusses multistream methods for combining multiple hidden Markov models to perform speech recognition. Finally, it introduces the ensemble methods used in the paper, bagging and boosting, in their basic form.

Section 3 discusses related work and their relation to our contributions, while Section 4 gives details about the data and the experimental protocols followed.

In the speech model considered, words are hidden Markov models composed of concatenations of phonetic hidden Markov models. In this setting it is possible to employ mixture models at any temporal level. Section 5 considers mixtures at the phoneme model level, where data with a phonetic segmentation is available. We can then restrict ourselves to a sequence classification problem in order to train a mixture model. Application of methods such as bagging and boosting to the phoneme classification task is then possible. However, using the resulting models for continuous speech recognition poses some difficulties in terms of complexity. Section 5.1 outlines how multistream decoding can be used to perform approximate inference in the resulting mixture model.

Section 6 discusses an algorithm, introduced in [3], for word error rate minimisation using boosting techniques.

While it appears trivial to do so by minimising some form of loss based on the word error rate, in practice successful application additionally requires use of a probabilistic model for inferring error probabilities in parts of misclassified sequences. The concepts of expected label and expected loss are introduced, of which the latter is used in place of the conventional loss. This integration of probabilistic models with boosting allows its use in problems where labels are not available.

Sections 7 and 8 conclude the paper with an extensive comparison between the proposed models. It is clearly shown neither of the boosting approaches employed manage to outperform a simple bagging model that is trained on presegmented phonetic data. Furthermore, in a follow-up experiment, we find that the performance of bagging when using triphone models achieves state-of-the-art results for the dataset used. These are significant findings, since most of the recent ensemble-based hidden Markov model research on speech recognition has focused invariably on boosting.

## 2. Background and Notation

Sequence learning and sequential decision making deal with the problem of modelling the relationship between sequential variables from a set of data and then using the models to make decisions. In this paper, we examine two types of sequence learning tasks: sequence classification and sequence recognition.

The sequence classification task entails assigning a sequence to one or more of a set of categories. More formally, we assume a finite label set  $\mathcal{Y}$  and a possibly uncountably infinite observation set  $\mathcal{X}$ . We denote the set of sequences of length  $n$  as  $\mathcal{X}^n \triangleq \times^n \mathcal{X}$  and the null sequence set by  $\mathcal{X}^0 \triangleq \emptyset$ . Finally, we denote the set of all sequences by  $\mathcal{X}^* \triangleq \bigcup_{n=0}^{\infty} \mathcal{X}^n$ . We observe sequences  $x = x_1, x_2, \dots$ , with  $x_i \in \mathcal{X}$  and  $x \in \mathcal{X}^*$ , and we use  $|x|$  to denote the length of a sequence  $x$ , while  $x_{t:T} = x_t, x_{t+1}, \dots, x_T$  denotes subsequences. In sequence classification, each  $x \in \mathcal{X}^*$  is associated with a label  $y \in \mathcal{Y}$ . A sequence classifier  $f \in \mathcal{F}$ , is a mapping  $f : \mathcal{X}^* \rightarrow \mathcal{Y}$ , such that  $f(x)$  corresponds to the predicted label, or classification decision, for the observed sequence  $x$ .

We focus on probabilistic classifiers, where the predicted label is derived from the *conditional* probability of the class given the observations, or *posterior class probability*  $\mathbf{P}(y | x)$ , with  $x \in \mathcal{X}^*$ ,  $y \in \mathcal{Y}$ , where we make no distinction between random variables and their realisations. More specifically, we consider a set of models  $\mathcal{M}$  and an associated set of observation densities and class probabilities  $\{p(x | y, \mu), \mathbf{P}(y | \mu) : \mu \in \mathcal{M}\}$  indexed by  $\mu$ . The posterior class probability according to model  $\mu$  can be obtained by using Bayes' theorem:

$$\mathbf{P}(y | x, \mu) = \frac{p(x | y, \mu) \mathbf{P}(y | \mu)}{p(x | \mu)}. \quad (1)$$

Any model  $\mu$  can be used to define a classification rule.

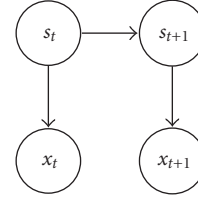


FIGURE 1: Graphical representation of a hidden Markov model, with arrows indicating dependencies between variables. The observations  $x_t$  and the next state  $s_{t+1}$  only depend on the current state  $s_t$ .

*Definition 1* (Bayes classifier). A classifier  $f_\mu : \mathcal{X}^* \rightarrow \mathcal{Y}$  that employs (1) and makes classification decisions according to

$$f_\mu(x) = \arg \max_{y \in \mathcal{Y}} \mathbf{P}(y | x, \mu) \quad (2)$$

is referred to as a Bayes classifier or a Bayes decision rule.

Formally, this task is exactly the same as nonsequential classification. The only practical difference is that the observations are sequences. However, care should be taken as this makes the implicit assumption that the costs of all incorrect decisions are equal.

In *sequence recognition*, we attempt to determine a sequence of events from a sequence of observations. More formally, we are given a sequence of observations  $x$  and are required to determine a sequence of labels  $y \in \mathcal{Y}^*$ , that is, the sequence  $y = y_1, y_2, \dots, y_k$ ,  $|y| \leq |x|$ , with maximum posterior probability  $\mathbf{P}(y | x)$ . In practice, models are used for which it is not necessary to exhaustively evaluate the set of possible label sequences. One such simple, yet natural, class is that of hidden Markov models.

### 2.1. Speech Recognition with Hidden Markov Models

*Definition 2* (hidden Markov model). A hidden Markov model (HMM) is a discrete-time stochastic process, with state variable  $s_t$  in some discrete space  $\mathcal{S}$ , and an observation variable  $x_t \in \mathcal{X}$ , such that

$$\begin{aligned} \mathbf{P}(s_t | s_{t-1}, s_{t-2}, \dots) &= \mathbf{P}(s_t | s_{t-1}), \\ \mathbf{P}(x_t | s_t, x_{t-1}, s_{t-1}, x_{t-2}, \dots) &= \mathbf{P}(x_t | s_t). \end{aligned} \quad (3)$$

The model is characterised by the observation distribution  $\mathbf{P}(x_t | s_t)$ , the transition distribution  $\mathbf{P}(s_t | s_{t-1})$ , and the initial state distribution  $\mathbf{P}(s_1) \equiv \mathbf{P}(s_1 | s_0)$ . These dependencies are shown graphically in Figure 1.

Training consists of two steps. First, select a class of hidden Markov models  $\mathcal{M}$ , with each model  $\mu \in \mathcal{M}$  corresponding to a pair of transition and observation densities  $\mathbf{P}(s_t | s_{t-1}, \mu)$ ,  $\mathbf{P}(x_t | s_t, \mu)$ . The second step is to select a model from  $\mathcal{M}$ . By additionally defining a prior density  $p(\mu)$  over  $\mathcal{M}$ ,

we can try to find the maximum *a posteriori* (MAP) model  $\mu^* \in \mathcal{M}$ , given a set of observation sequences  $\mathcal{D}$

$$\mu^* = \arg \max_{\mu \in \mathcal{M}} p(\mu | \mathcal{D}). \quad (4)$$

The class  $\mathcal{M}$  is restricted to models with a particular number of states and allowed transitions between states. In this paper, the optimisation is performed through expectation maximisation.

The most common way to apply such models to speech recognition is to associate each state  $s$  with phonological units  $a \in \mathcal{A}$ , such as phonemes, syllables, or words, through a distribution  $\mathbf{P}(a | s)$ , which takes values in  $\{0, 1\}$  in usual practice; thus, each state is mapped to only one phoneme. This is done by modelling each phoneme as a small HMM (Figure 2) and combining them into a larger HMM, such as the one shown in Figure 3, with a set of parallel chains such that each chain maps to one word; for example, given that we are in the state  $s = 4$  at some time  $t$ , then we are also definitely (i.e., with probability 1) in Word A and Phoneme B at time  $t$ . In general, if we can determine the probabilities for sequences of states, we can also determine the most probable sequence of words or phonemes; that is, given a sequence of observations  $x_{1:T}$ , we calculate the state distribution  $\mathbf{P}(s_{1:T} | x_{1:T})$  and subsequently a distribution over phonologies, to wit the probabilities of possible word, syllable, or phoneme sequences. Thus, the problem of recognising word sequences is reduced to the problem of state estimation.

**2.2. Multistream Decoding.** When we wish to combine evidence from  $n$  different models, state estimation is significantly harder, as the number of effective states is  $|\mathcal{S}|^n$ . However, multistream decoding techniques can be used as an approximation to the full mixture model [4]. Such techniques derive their name from the fact that they were originally used to combine models which had been trained on different streams of data or features [5]. In this paper, we instead wish to combine evidence from models trained on different samples of the same data.

In multistream decoding each subunit model corresponding to a phonological unit  $a$  is comprised of  $n$  submodels  $a = \{a_i : i \in [1, n]\}$  associated with the subunit level at which the recombination of the input streams should be performed. For any given  $a$  and a distribution over models  $\pi(a_i | a)$ , the observation density conditioned on the unit  $a$  can be written as

$$\pi(x | a) = \sum_{i=1}^n p(x | a_i) \pi(a_i | a), \quad (5)$$

where  $\pi(a_i | a)$  can be seen as a weight for expert  $i$ . This may vary across  $a$ , but herein we consider the case where the weight is fixed, that is,  $\pi(a_i | a) = w_i$  for all  $a$ . We consider *state-locked* multistream decoding, where all submodels are forced to be at the same state. This can be viewed as creating another Markov model with emission distribution

$$\pi(x_t | s_t, a) = \sum_{i=1}^n p(x_t | s_t, a_i) \pi(a_i | a). \quad (6)$$

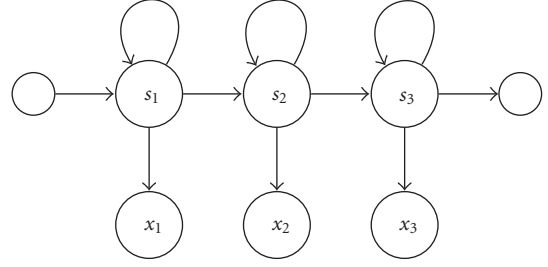


FIGURE 2: Graphical representation of a phoneme model with 3 emitting states, as well as initial and terminal nonemitting states. The arrows depict dependencies between specific states. All the phoneme models used in this paper employed the above topology.

An alternative is the exponentially weighted product of emission distributions:

$$\pi(x_t | s_t, a) = \prod_{i=1}^n p(x_t | s_t, a_i)^{\pi(a_i | a)}. \quad (7)$$

However, this approximation does not arise from (5) but from assuming a factorisation of the observations  $p(x_t | s_t) = \prod_{i=1}^n p(x_t^i | s_t)$ , which is useful when there is a different model for different parts of the observation vector.

Multistream techniques are hardly limited to the above. For example, Misra et al. [6] describe a system where  $\pi$  is related to the entropy of each submodel, while Ketabdar et al. [7] describe a multistream method utilising state posteriors. We, however, shall concentrate on the two techniques outlined above, as well as a single-stream technique to be described in Section 5.1.

**2.3. Ensemble Methods.** We investigate the use of ensemble methods in the class of *static* mixture models for speech recognition. Such methods construct an aggregate model from a set of base hypotheses  $M \triangleq \{\mu_i : i = 1, \dots, N\}$ . Each hypothesis  $\mu_i$  indexes a set of conditional distributions  $\{\mathbf{P}(\cdot | \cdot, \mu_i) : i = 1, \dots, N\}$ . To complete the model, we employ a set of weights  $W \triangleq \{w_i : i = 1, \dots, N\}$  corresponding to the probability of each base hypothesis, so that  $w_i \triangleq \mathbf{P}(\mu_i)$ . Thus, we can form a mixture model, assuming  $\mathbf{P}(\mu_i | x) = \mathbf{P}(\mu_i)$  for all  $x \in \mathcal{X}^*$ :

$$\mathbf{P}(\cdot | \cdot, M, W) = \sum_{i=1}^N w_i \mathbf{P}(\cdot | \cdot, \mu_i). \quad (8)$$

Two questions that arise when training such models are how to select  $M$  and  $W$ . In this paper, we consider two different approaches, bagging and boosting.

**2.3.1. Bagging.** Bagging [8] can be seen as a method for sampling the model space  $\mathcal{M}$ . We first require a learning algorithm  $\Lambda : (\mathcal{X}^* \times \mathcal{Y})^* \rightarrow \mathcal{M}$  that maps (While we restrict ourselves to the deterministic case for simplicity, bagging is applicable to stochastic learning algorithms as well.) from a *dataset*  $D \in (\mathcal{X}^* \times \mathcal{Y})^*$  of data pairs  $(x, y)$

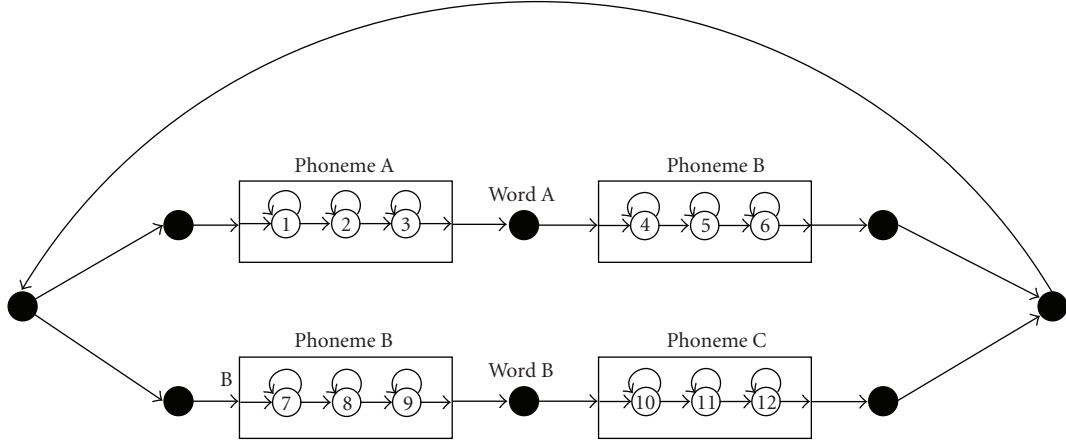


FIGURE 3: A hidden Markov model for speech recognition. The figure depicts how models of three phonemes, A, B, C, are used to construct a single hidden Markov model for distinguishing between two different words. The states are indexed uniquely. Black circles indicate non-emitting states.

to models  $\mu \in \mathcal{M}$ . We then sample  $N$  datasets  $D_i$  from a distribution  $\mathcal{D}$ , for  $i = 1, \dots, N$ . For each  $D_i$ , the learning algorithm  $\Lambda$  generates a model  $\mu_i \triangleq \Lambda(D_i)$ . The models  $M \triangleq \{\mu_i : i = 1, \dots, N\}$  can be combined into a mixture with  $w_i = 1/N$  for all  $i$ :

$$\mathbf{P}(y | x, M, W) = \frac{1}{N} \sum_{i=1}^N \mathbf{P}(y | x, \mu_i). \quad (9)$$

In Bagging,  $D_i$  is generated by sampling with replacement from the original dataset  $D$ , with  $|D_i| = |D|$ . Thus,  $D_i$  is a *bootstrap replicate* of  $D$ .

**2.3.2. Boosting.** Boosting algorithms [9–11] are another family of ensemble methods. The most commonly used boosting algorithm for classification is AdaBoost [9]. Though many variants of AdaBoost for multiclass classification problems exist, in this paper we will use AdaBoost.M1.

An AdaBoost ensemble is a mixture model composed of  $N$  models  $\mu_i$  and weights  $w_i$ , as in the previous section. The models and weights are created in an iterative manner. At iteration  $j$ , the model  $\mu_j \triangleq \Lambda(D_j)$  is created from a *weighted* bootstrap sample  $D_j$  of the training dataset  $D = \{d_i : i \in [1, n]\}$ , with  $d_i = (x_i, y_i)$ . The probability of adding example  $d_i$  to the bootstrap replicate  $D_j$  is denoted as  $p_j(d_i)$ , with  $\sum_i p_j(d_i) = 1$ . At the end of iteration  $j$  of AdaBoost.M1,  $\beta_j$  is calculated according to

$$\beta_j = \ln \frac{1 - \varepsilon_j}{\varepsilon_j}, \quad (10)$$

where  $\varepsilon_j \triangleq \sum_i p_j(d_i) \ell(d_i)$  is the empirical expected loss of the  $j$ th, with  $\ell(d_i) \triangleq \mathbb{I}\{h_i \neq y_i\}$  being the *sample loss* of example  $d_i$ , where  $\mathbb{I}\{\cdot\}$  is an indicator function. At the end of each iteration, sampling probabilities are updated according to

$$p_{j+1}(d_i) = \frac{p_j(d_i) \exp(\beta_j \ell(d_i))}{Z_j}, \quad (11)$$

where  $Z_j \triangleq \sum_i p_j(d_i) \exp(\beta_j \ell(d_i))$  is a normalisation factor. Thus, incorrectly classified examples are more likely to be included in the next bootstrap data set. The final model is a mixture with  $N$  components  $\mu_i$  and weights  $w_i \triangleq \beta_i / \sum_{j=1}^N \beta_j$ .

### 3. Contributions and Related Work

The original AdaBoost algorithm had been defined for classification and regression tasks, with the regression case receiving more attention recently (see [10] for an overview). In addition, research in the application of boosting to sequence learning and speech recognition has intensified [12–15]. The application of other ensemble methods, however, has been limited to random decision trees [16, 17]. In our view, bagging [8] is a method that has been somewhat unfairly neglected, and we present results that show that it can outperform boosting in an unbiased experiment.

One of the simplest ways to apply ensemble methods to speech recognition is to employ them at the state level. For example, Schwenk [18] proposed a HMM/artificial neural network (ANN) system, with the ANNs used to compute the posterior phoneme probabilities at each state. Boosting itself was performed at the ANN level, using AdaBoost with confidence-rated predictions, using the frame error rate as the sample loss function. The resulting decoder system differed from a normal HMM/ANN hybrid in that each ANN was replaced by a mixture of ANNs that had been provided via boosting. Thus, such a technique avoids the difficulties of performing inference on mixtures, since the mixtures only model instantaneous distributions. Zweig and Padmanabhan [19] appear to be using a similar technique, based on Gaussian mixtures. The authors additionally describe a few boosting variants for large-scale systems with thousands of phonetic units. Both papers report mild improvements in recognition.

One of the first approaches to utterance-level boosting is due to Cook and Robinson [20], who employed a boosting scheme, where the sentences with the highest error rate were



classified as “incorrect” and the rest “correct,” irrespective of the absolute word error rate of each sentences. The weights of all frames constituting a sentence were adjusted equally and boosting was applied at the frame level. This however does not manage to produce as good results as the other schemes described by the authors. In our view, which is partially supported by the experimental results in Section 6, this could have been partially due to the lack of a temporal credit assignment mechanism such as the one we present. An early example of a nonboosting approach for the reduction of word error rate is [21], which employed a “corrective training scheme.”

In related work on utterance-level boosting, Zhang and Rudnicky [22] compared use of the posterior probability of each possible utterance for adjusting the weights of each utterance with a “nonboosting” method, where the same weights are adjusted according to some function of the word error rate. In either case, utterance posterior probabilities are used for recombining the experts. Since the number of possible utterances is very large, not all possible utterances are used but an  $N$ -best list. For recombination, the authors consider two methods: firstly, choosing the utterance with maximal sum of weighted posterior (where the weights have been determined by boosting). Secondly, they consider combining via ROVER, a dynamic programming method for combining multiple speech recognisers (see [23]). Since the authors’ use of ROVER entails using just one hypothesis from each expert to perform the combination, in [15] they consider a scheme where the  $N$ -best hypotheses are reordered according to their *estimated* word error rate. In further work [24] the authors consider a boosting scheme for assigning weights to frames, rather than just to complete sentences. More specifically, they use the currently estimated model to obtain the probability that the correct word has been decoded at any particular time, that is, the posterior probability that the word at time  $t$  is  $a_t$  given the model and the sequence of observations. In our case we use a slightly different formalism in that we calculate the expectation of the loss according to an independent model.

Finally, Meyer and Schramm [13] propose an interesting boosting scheme with a weighted sum model recombination. More precisely, the authors employ AdaBoost.M2 at the utterance level, utilising the posterior probability of each utterance for the loss function. Since the algorithm requires calculating the posterior of every possible class (in this case an utterance) given the data, exact calculation is prohibitive. The required calculation however can be approximated by calculating the posterior only for the subset of the top  $N$  utterances and assuming the rest are zero. Their model recombination scheme relies upon treating each expert as a different pronunciation model. This results in essentially a mixture model in the form of (5), where the weight of each expert is derived from the boosting algorithm. They further robustify their approach through a language model. Their results indicate a slight improvement (in the order of 0.5%) in a large vocabulary continuous speech recognition experiment.

More recently, an entirely different and interesting class of complementary models were proposed in [12, 16, 17].

The core idea is the use of randomised decision trees to create multiple experts, which allows for more detailed modelling of the strengths and weaknesses of each expert, while [12] presents an extensive array of methods for recombination during speech recognition. Other recent work has focused on slightly different applications. For example, a boosting approach for language identification was used in [14, 25], which utilised an ensemble of Gaussian mixture models for both the *target class* and the *antimodel*. In general, however, bagging methods, though mentioned in the literature, do not appear to be used, and recent surveys, such as [12, 26, 27] do not include discussions of bagging.

**3.1. Our Contribution.** This paper presents methods and results for the use of both boosting and bagging for phoneme classification and speech recognition. Apart from synthesising and extending our previous results [2, 3], the main purpose of this paper is to present an *unbiased* experimental comparison between a large number of methods, controlling for the appropriate choice of hyperparameters and using a principled statistical methodology for the evaluation of the significance of the results. If this is not done, then it is possible to draw incorrect conclusions.

Section 5 describes our approach for phoneme-level training of ensemble methods (boosting and bagging). In the phoneme classification case, the formulation of the task is essentially the same as that of static classification; the only difference is that the observations are sequences rather than single values. As far as we know, our past work [2] is the only one employing ensemble methods at the phoneme level. In Section 5, we extend our previous results by comparing boosting and bagging in terms of both classification and recognition performance and show, interestingly, that bagging achieves the same reduction in recognition error rates as boosting, even though it cannot match boosting classification error rate reduction. In addition, the section compares a number of different multistream decoding techniques.

Another interesting way to apply boosting is to use it at the sentence level, for the purposes of explicitly minimising the word error rate. Section 6 presents a boosting-based approach to minimise the word error rate originally introduced in [3].

Finally, Section 7 presents an extensive, unbiased experimental comparison, with separate model selection and model testing phase, between the proposed methods and a number of baseline systems. This shows that the simple phoneme-level bagging scheme outperforms all of the other boosting schemes explored in this paper significantly. Finally, further results using tri-phone models indicate that state-of-the-art performance is achievable for this dataset using bagging but not boosting.

## 4. Data and Methods

The phoneme data was based on a presegmented version of the OGI Numbers 95 (N95) data set [28]. This data set was converted from the original raw audio data into a set

of features based on Mel-Frequency Cepstrum Coefficients (MFCC) [29] (with 39 components, consisting of three groups of 13 coefficients, namely, the static coefficients and their first and second derivatives) that were extracted from each frame. The data contains 27 distinct phonemes (or 80 tri-phones in the tri-phone version of the dataset) that compose 30 dictionary words. There are 3233 training utterances and 1206 test utterances, containing 12510 and 4670 words, respectively. The segmentation of the utterances into their constituent phonemes resulted in 35562 training segments and 12613 test segments, totalling 486537 training frames and 180349 test frames, respectively. The feature extraction and phonetic labelling are described in more detail in [30].

**4.1. Performance Measures.** The comparative performance measure used depends on the task. For the phoneme classification task, the *classification error* is used, which is the percentage of misclassified examples in the training or testing data set. For the speech recognition task, the following *word error rate* is used:

$$\text{WER} = \frac{N_{\text{ins}} + N_{\text{sub}} + N_{\text{del}}}{N_{\text{words}}}, \quad (12)$$

where  $N_{\text{ins}}$  is the number of word insertions,  $N_{\text{sub}}$  the number of word substitutions, and  $N_{\text{del}}$  the number of word deletions. These numbers are determined by finding the minimum number of insertions, substitutions, or deletions necessary to transform the target utterance into the emitted utterance for each example and then summing them for all the examples in the set.

**4.2. Bootstrap Estimate for Speech Recognition.** In order to establish the significance of the reported results, we employ a bootstrap method; (see [31]). More specifically, we use the approach suggested by Bisani and Ney [32] for speech recognition. It amounts to using the results of speech recognition on a test set of sentences as an empirical distribution of errors. Using this method, we obtain a bootstrap estimate of the probability distribution of the difference in word error rate  $\Delta W$  between two systems, from  $B$  bootstrap samples  $\Delta W_k$  of the word error rate difference:

$$\mathbf{P}(\Delta W > u) = \int_u^\infty p(\Delta W) d\Delta W \approx \frac{1}{B} \sum_{k=1}^B \mathbb{I}\{\Delta W_k > u\}, \quad (13)$$

where  $\mathbb{I}\{\cdot\}$  is an indicator function. This approximates the probability that system  $A$  is better than system  $B$  by more than  $u$ . See [31] for more on the properties of the bootstrap and [33] for the convergence of empirical processes and their relation to the bootstrap.

**4.3. Parameter Selection.** The models employed have a number of hyperparameters. In order to perform unbiased comparisons, we split the training data into a smaller training set of 2000 utterances and a hold-out set of 1233 utterances. For the preliminary experiments performed in Sections 5 and 6, we train all models on the small training set and report the performance on both the training and the hold-out set. For the experiments in Section 7, each model's

hyperparameters are selected independently on the hold-out set. Then the model is trained on the complete training set and evaluated in the independent test set.

For the classification task (Section 5), we used presegmented data. Thus, the classification could be performed using a Bayes classifier composed of 27 hidden Markov models, each one corresponding to one class. Each phonetic HMM was composed of the same number of hidden states (And an additional two nonemitting states: the initial and final states.), in a left-to-right topology, and the distributions corresponding to each state were modelled with a Gaussian mixture model, with each Gaussian having a diagonal covariance matrix. In Section 5.2, we select the number of states per phoneme from  $\{1, 2, 3, 4, 5\}$  and the mixture components from  $\{10, 20, 30, 40\}$  in the hold-out set for a single HMM and then examine whether bagging or boosting can improve the classification or speech recognition performance.

In all cases, the diagonal covariance matrix elements of each Gaussian were clamped to a lower limit of 0.2 times the global variance of the data. For continuous speech recognition, transitions between word models incurred an additional likelihood penalty of  $\exp(-15)$  while calculating the most likely sequence of states. Finally, in all continuous speech recognition tasks, state sequences were constrained to remain in the same phoneme for at least three acoustic frames.

For phoneme-level training, the adaptation of each phoneme model was performed in two steps. Firstly, the acoustic frames belonging to each phonetic segment were split into a number of equally sized intervals, where the number of intervals was equal to the number of states in the phonetic model. The Gaussian mixture components corresponding to the data for each interval were initialised via 25 iterations of the K-means algorithm (see, e.g., [34]). After this initialisation was performed, a maximum of 25 iterations of the EM algorithm were run on each model, with optimisation stopping earlier if, at any point in time  $t$ , the likelihood  $L_t$  satisfied the stopping criterion  $(L_t - L_{t-1})/L_t < \epsilon$ , with  $\epsilon = 10^{-5}$  being used in all experiments that employed EM for optimisation.

For the utterance-level training described in Section 6, the same initialisation was performed. The inference of the final model was done through expectation maximisation (using the Viterbi approximation) on concatenated phonetic models representing utterances. Note that performing the full EM computation is costlier and does not result in significantly better generalisation performance, at least in this case. The stopping criterion and maximum iterations were the same as those used for phoneme-level training.

Finally, the results in Section 7 present an unbiased comparison between models. In order to do this, we selected the parameters of each model, such as the number of Gaussians and number of experts, using the performance in the hold-out set. We then used the selected parameters to train a model on the full training dataset. The models were evaluated on the separate testing dataset and compared using the bootstrap estimate described in Section 4.2.

## 5. Phoneme-Level Bagging and Boosting

A simple way to apply ensemble techniques such as bagging and boosting is to cast the problem into the classification framework. This is possible at the phoneme level, where each class  $y \in \mathcal{Y}$  corresponds to a phoneme. As long as the available data are annotated so that subsequences containing single phoneme data can be extracted, it is natural to adapt each hidden Markov model  $\mu_y$  to a single class  $y$  out of the possible  $|\mathcal{Y}|$ , where  $|\cdot|$  denotes the cardinality of the set, and combine the models into a Bayes classifier in the manner described in Section 2. Such a Bayes classifier can then be used as an expert in an ensemble.

In both cases, each example  $d$  in the training dataset  $D$  is a sequence segment corresponding to data from a single phoneme. Consequently, each example  $d$  has the form  $d = (x, y)$ , with  $x \in \mathcal{X}^*$  being a subsequence of features corresponding to single phoneme data and  $y \in \mathcal{Y}$  being a phoneme label.

Both methods iteratively construct an ensemble of  $N$  models. At each iteration  $j$ , a new classifier  $h_j$  is created, consisting of a set of hidden Markov models:  $h_j = \{\mu_1^j, \mu_2^j, \dots, \mu_{|\mathcal{Y}|}^j\}$ . Each model  $\mu_y^j$  is adapted to the set of examples  $\{d_k \in \mathcal{D}_j \mid y_k = y\}$ , where  $\mathcal{D}_j$  is a bootstrap replicate of  $\mathcal{D}$ . In order to make decisions, the experts are weighted by the mixture coefficients  $w_i \triangleq \pi(h_i)$ . The only difference between the two methods is the distribution that  $\mathcal{D}_j$  is sampled from and the definition of the coefficients.

For “*bagging*”,  $\mathcal{D}_j$  is sampled *uniformly* from  $\mathcal{D}$ , and the probability over the mixture components is also uniform, that is,  $\pi(h_i) = N^{-1}$ .

For “*boosting*”,  $\mathcal{D}_j$  is sampled from  $\mathcal{D}$  using the distribution defined in (11), while the expert weights are defined as  $\pi(h_i) = \beta_i / \sum_j \beta_j$ , where  $\beta$  is given by (10). The AdaBoost method used was AdaBoost.M1.

Since previous studies in nonsequential classification problems had shown that an increase in generalisation performance may be obtained through the use of those two ensemble methods, it was expected that they would have a similar effect on performance in phoneme classification tasks. This is tested in Section 5.2. While using the resulting phoneme classification models for continuous speech recognition is not straightforward, we describe some techniques for combining the ensembles resulting from this training in order to perform sequence recognition in Section 5.1.

**5.1. Continuous Speech Recognition with Mixtures.** The approach described is easily suitable for phoneme classification, since each phonetic model is now a mixture model (Figure 4), which can be used to classify phonemes given presegmented data. However, the phoneme mixtures can also be combined into a speech recognition mixture. Thus, we can still employ ensemble methods for the full speech recognition problem by training with segmented data to produce a number of expert models which can then be recombined during decoding on unsegmented data.

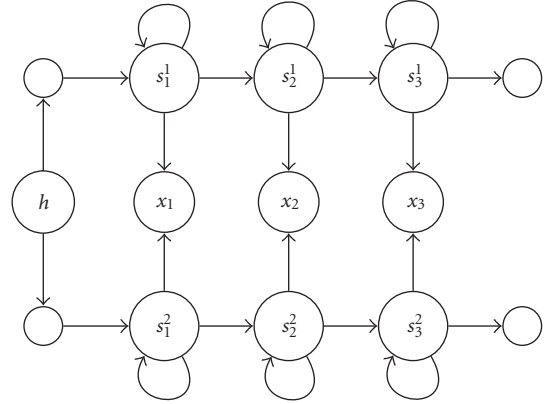


FIGURE 4: A phoneme mixture model. The generating model depends on the hidden variable  $h$ , which determines the mixing coefficients between model 1 and 2. The random variable  $h$  may in general depend on other variables. The distribution of the observation is a mixture between the two distributions predicted by the two hidden models, mixed according to the mixture model  $h$ .

The first technique employed for sequence decoding uses an HMM comprising all phoneme models created during the boosting process, connected in the manner shown in Figure 5. Each phase of the boosting process creates a sub-model  $i$ , which we will refer to as *expert* for disambiguation purposes. Each expert is a classification model that employs one hidden Markov model for each phoneme. For some sequence of observations, each expert calculates the posterior probability of each phonetic class given the observation and its model. Two types of techniques are considered for employing the models for inferring a sequence of words.

In the *single-stream* case, decoding is performed using the Viterbi algorithm in order to find a sequence of states maximising the posterior probability of the sequence. A normal hidden Markov model is constructed in the way shown in Figure 5, with each phoneme being modelled as a mixture of expert models. In this case we are trying to find the sequence of states  $\{s_t = s_t^j\}$  with maximum likelihood. The transition probabilities leading from anchor states (black circles in the figure) to each model are set to  $w_i = \pi(h_i)$ .

This type of decoding would have been appropriate if the original mixture had been inferred as a type of switching model, where only one submodel is responsible for generating the data at each point in time and where switching between models can occur at anchor states.

The models may also be combined using *multistream* decoding (see Section 2.2). The advantage of such a method is that it uses information from all models. The disadvantage is that there are simply too many states to be considered. In order to simplify this, we consider multistream decoding synchronised at the state level, that is, with the constraint that  $P(s_t^i \neq s_t^j) = 0$  if  $j \neq i$ . This corresponds to (5), where the weight of stream  $i$  is again  $w_i$ .

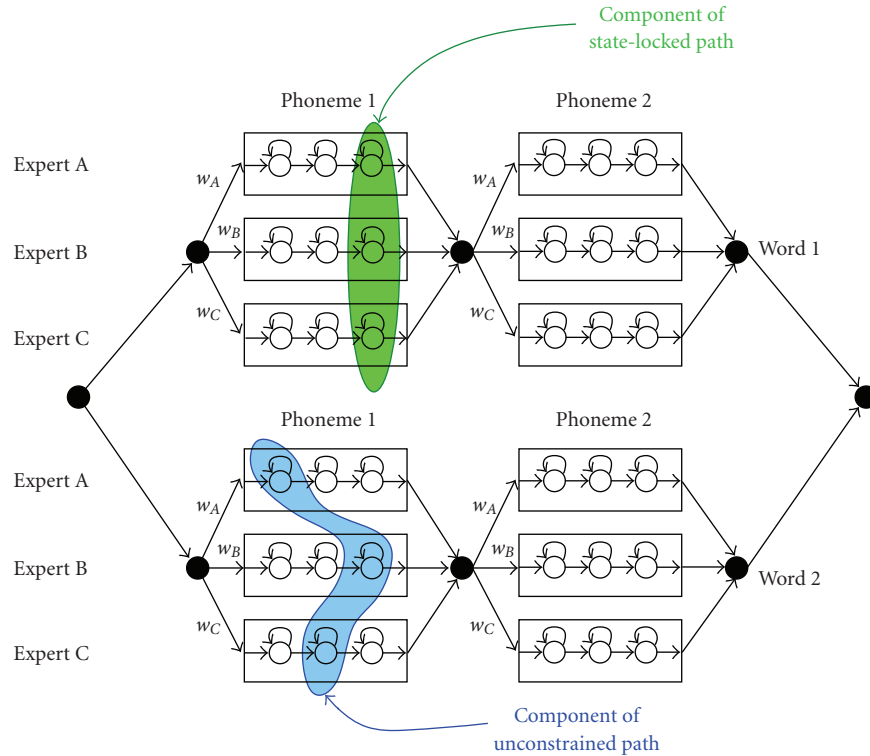


FIGURE 5: Single-path multistream decoding for two vocabulary words consisting of two phonemes each. When there is only one expert, the decoding process is done normally. In the multiple-expert case, phoneme models from each expert are connected in parallel. The transition probabilities leading from the anchor states to the hidden Markov model corresponding to each experts are the weights  $w_i$  of each expert.

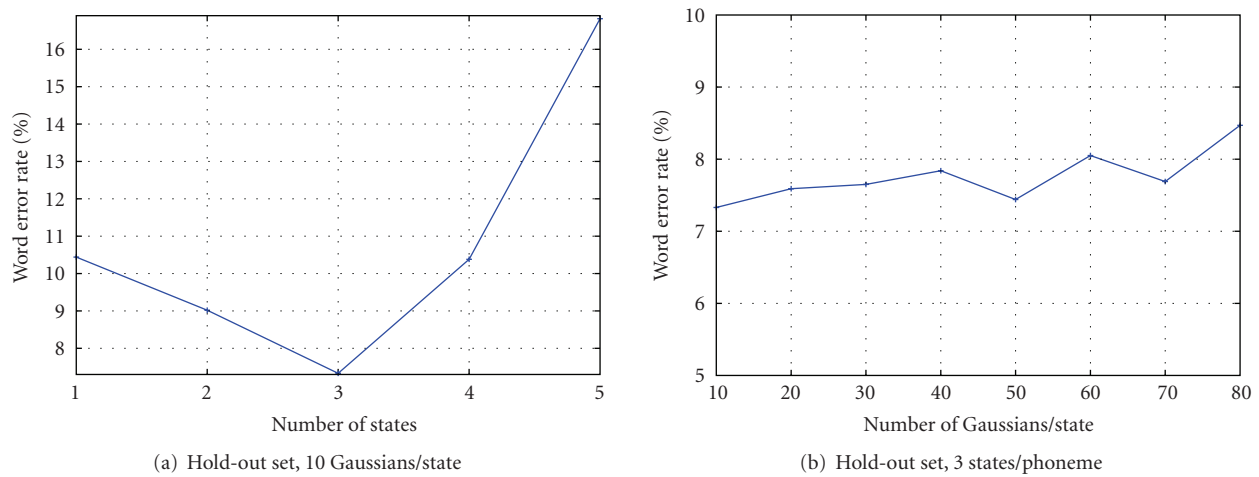


FIGURE 6: In the experiments reported in Section 5.2, the number of states and number of Gaussian mixtures per state were tuned on a hold-out set prior to the analysis. (a) displays the word error rate performance of an HMM with 10 Gaussians per state when the number of emitting states per phoneme is varied, with rather dramatic effects. (b) displays the word error rate performance of an HMM with 3 emitting states as the number of Gaussians per state varies. In this case, the effect on generalisation is markedly lower.

**5.2. Experiments with Boosting and Bagging Phoneme-Level Models.** The experiments described in this section were performed with a fixed number of states for all phonemes, as well as with a fixed number of Gaussians per state. The selection of these hyperparameters was performed on a hold-out set, as described in Section 4. The hold-out

set results are shown in Figure 6. After selecting those hyperparameters, we perform an exploratory comparison (An experiment that uses an unbiased procedure to select the number of experts independently for boosting and bagging is described in Section 7.) of the performance of boosting and bagging as the number of mixture components are increased,



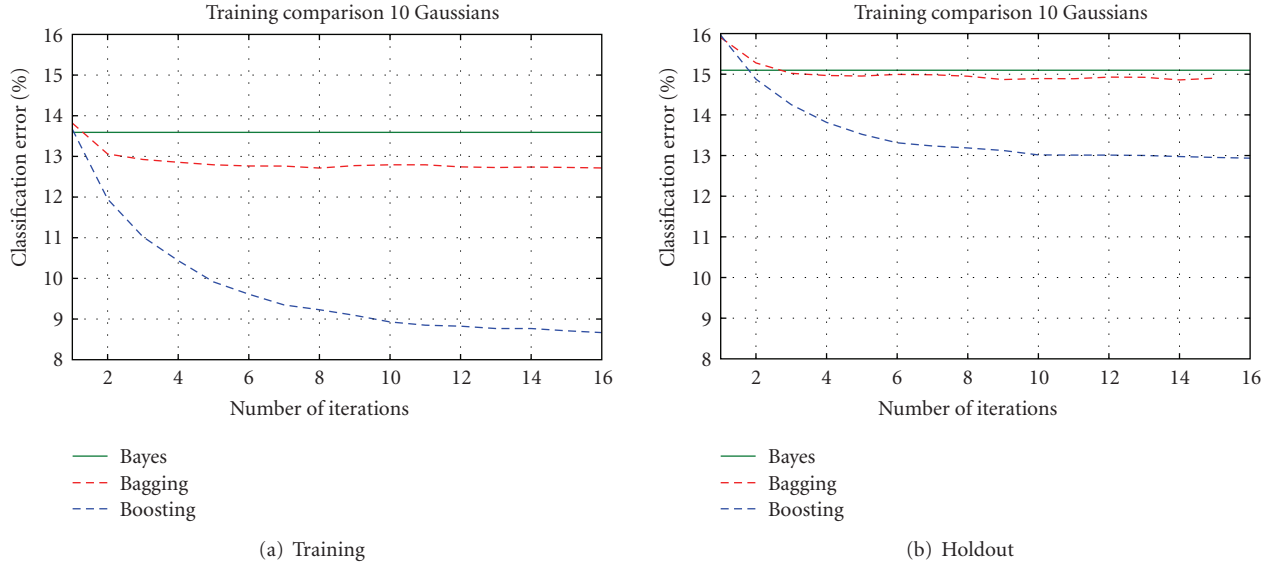


FIGURE 7: Classification errors for a bagged and a boosted ensemble of Bayes Classifiers as the number of experts is increased. For reference, the corresponding errors for a single Bayes Classifier trained on the complete training set are also included. There were 10 Gaussians per state and 3 states per phoneme for all models.

for the tasks of phoneme classification and speech recognition. For the latter problem, we also examine the relative merits of different decoding techniques.

Since the available data includes segmentation information, it makes sense to first limit the task to training for phoneme classification. This enables the direct application of ensemble training algorithms by simply using each segment as a training example.

Two methods were examined for this task: bagging and boosting. At each iteration of either method, a sample from the training set was made according to the distribution defined by either algorithm and then a Bayes classifier composed of  $|\mathcal{Y}|$  hidden Markov models, one for each phonetic class  $y \in \mathcal{Y}$ , was trained.

It then becomes possible to apply the boosting and bagging algorithms by using Bayes Classifiers as the experts. The N95 data was presegmented into training examples, so that each one was a segment containing a single phoneme. Bootstrapping was performed by sampling through these examples. The classification error of each classifier was used to calculate the boosting weights. The test data was also segmented in subsequences consisting of single phoneme data, so that the models could be tested on the phoneme classification tasks.

Figure 7 compares the classification performance of bagging and boosting as the number of experts increases with that of the Bayes classifier trained on the full training data. As can be seen in Figure 7(a), both bagging and boosting manage to reduce the phoneme classification error considerably in the training, with boosting continuing to make improvements until the maximum number of iterations. For bagging, the improvement in classification was limited to the first 4 iterations, after which performance remained

constant. The situation was similar when comparing the models in the hold-out set, shown in Figure 7(b). There, however, bagging failed to improve upon the baseline system.

Finally, an exploratory comparison between the models on the task of continuous speech recognition was made. This was necessary, in order to decide on a method for performing decoding when dealing with multiple models. The three relatively simple methods of single-stream and multistream decoding (the latter employing either weighted product or weighted sum) were evaluated on the hold-out set. As can be seen in Figure 8, the weighted sum method consistently performed the best for both bagging and boosting. This was expected since it was the only method with some justification in our particular case, as it arises out of constraining the full state inference problem on the mixture. The multistream product method is not justified here, since each model had exactly the same observation variables. The single-stream model could perhaps be justified under the assumption of a switching model, where a different expert can be responsible for the observations in each phoneme. That might explain the fact that its performance is not degrading in the case of bagging, as the components of each mixture should be quite similar to each other, something which is definitely not the case with boosting, where each model is trained on a different distribution of the data.

A fuller comparison between bagging and boosting at the phoneme level will be given in Section 7, where the number of Gaussian units per state and the number of experts will be independently tuned on the hold-out set and evaluated on a separate test set. There, it will be seen that with an *unbiased* hyperparameter selection, bagging actually outperforms boosting.

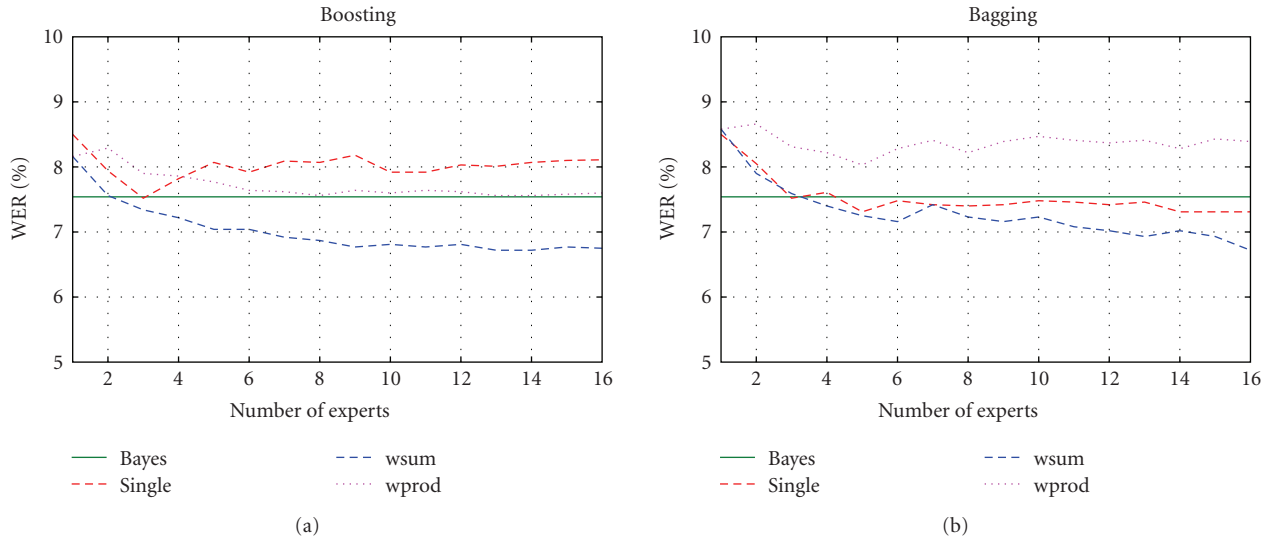


FIGURE 8: Generalisation performance on the hold-out set in terms of word error rate after training with segmentation information. Results are shown for both boosting and bagging, using three different methods for decoding. Single-path and multistream. Results are shown for three different methods single-stream (*single*), and state-locked multistream using either a weighted product (*wprod*) or weighted sum (*wsum*) combination.

## 6. Expectation Boosting for WER Minimisation

It is also possible to apply ensemble training techniques at the utterance level. As before, the basic models used are HMMs that employ Gaussian mixtures to represent the state observation distributions. Attention is restricted to boosting algorithms in this case. In particular, we shall develop a method that uses boosting to simultaneously utilise information about the complete utterance, together with an estimate about the phonetic segmentation. Since this estimate will be derived from bootstrapping our own model, it is unreliable. The method developed will take into account this uncertainty.

More specifically, similarly to [20], sentence-level labels (sequences of words without time indications) are used to define the error measure that we wish to minimise. The measure used is related to the word error rate, as defined in (12). In addition to a loss function at the sentence level, a probabilistic model is used to define a distribution for the loss at the frame level. Combined, the two can be used for the greedy selection of the next base hypothesis. This is further discussed in the following section.

**6.1. Boosting for Word Error Rate Minimisation.** In the previous section (and [2]) we have applied boosting to speech recognition at the phoneme level. In that framework, the aim was to reduce the *phoneme classification error* in presegmented examples. The resulting boosted phoneme models were combined into a single speech recognition model using *multistream* techniques. It was hoped that we could reduce the word error rate as a side effect of performing better phoneme classification, and three different approaches were examined for combining the models in order to perform continuous speech recognition. However,

since the measure that we are trying to improve is the word error rate and since we did not want to rely on the existence of segmentation information, minimising the word error rate directly would be desirable. This section describes such a scheme using boosting techniques.

We describe a training method that we introduced in [3], specific to boosting and hidden Markov models (HMMs), for word error rate reduction. We employ a score that is exponentially related to the word error rate of a sentence example. The weights of the frames constituting a sentence are adjusted depending on our expectation of how much they contribute to the error. Finally, boosting is applied at the sentence and frame level simultaneously. This method has arisen from a twofold consideration: firstly, we need to directly minimise the relevant measure of performance, which is the word error rate. Secondly, we need a way to more exactly specify which parts of an example most probably have contributed to errors in the final decision. Using boosting, it is possible to focus training on parts of the data which are most likely to give rise to errors while at the same time doing it in such a manner as take into account the actual performance measure. We find that both aspects of training have an important effect.

Section 6.1.1 describes word error rate-related loss functions that can be used for boosting. Section 6.1.2 introduces the concept of *expected error*, for the case when no labels are given for the examples. This is important for the task of word error rate minimisation. Previous sections on HMMs and multistream decoding described how the boosted models are combined for performing the speech recognition task. Experimental results are outlined in Section 6.2. We conclude with an experimental comparison between different methods in Section 7, followed by a discussion.

**6.1.1. Sentence Loss Function.** A commonly used measure of optimality for speech recognition tasks is the word error rate (12). We would like to minimise this quantity using boosting techniques. In order to do this, a dataset is considered, where each example is a complete sentence and where the loss  $\ell(d)$  for each example  $d$  is given by some function of the word error rate for the sentence.

The word error rate for any particular sentence can take values in  $[0, \infty)$ , while the AdaBoost algorithm that is employed herein requires a sample loss function with range  $[-1, 1]$ . For this reason we employ the *ad hoc*, but reasonable, mapping  $\ell : [0, \infty) \rightarrow [-1, 1]$

$$\ell(x) = 1 - 2e^{-\eta x}, \quad (14)$$

where  $x$  is the word error rate. When  $\ell(x) = -1$ , an example is considered as classified correctly, and, when  $\ell(x) = 1$ , the example is considered to be classified incorrectly. This mapping includes a free parameter  $\eta > 0$ . Increasing the parameter  $\eta$  increases the sharpness of the transition, as shown in Figure 9. This function is used for  $\ell(\cdot)$  in (11).

While this scheme may well result in some improvement in word recognition with boosting, while avoiding relying on potentially erroneous phonetic labels, there is some information that is not utilised. Knowledge of the required sequence of words, together with the obtained sequence of words for each decoded sentence results in a set of errors that are fairly localised in time. The following sections discuss how it is possible to use a model that capitalises on such knowledge in order to define a distribution of errors over time.

**6.1.2. Error Expectation for Boosting.** In traditional supervised settings we are provided with a set of examples and labels, which constitute our training set, and thus it is possible to apply algorithms such as boosting. However, this becomes problematic when labels are noisy; (see, e.g., [35]). Such an example is a typical speech recognition data set. Most of the time such a data set is composed of a set of sentences, with a corresponding set of transcriptions. However, while the transcriptions may be accurate as far as the intention of the speakers or the hearing of the transcriber is concerned, subsequent translation of the transcription into phonetic labels is bound to be error prone, as it is quite possible for either the speaker to mispronounce words, or for the model that performs the automatic segmentation to make mistakes. In such a situation, adapting a model so that it minimises the errors made on the segmented transcriptions might not automatically lead into a model that minimises the word error rate, which is the real goal of a speech recognition system.

For this purpose, the concept of error expectation is introduced. Thus, rather than declaring with absolute certainty that an example is incorrect or not, we simply define  $\ell(d_i) = \mathbf{P}(y_i \neq h_i)$ , so that the sample loss is now the probability that a mistake was made on example  $i$  and we consider  $y_i$  to be a random variable. Since boosting can admit any sample loss function [9], this is perfectly reasonable, and it is possible to use this loss as a sample loss in a boosting

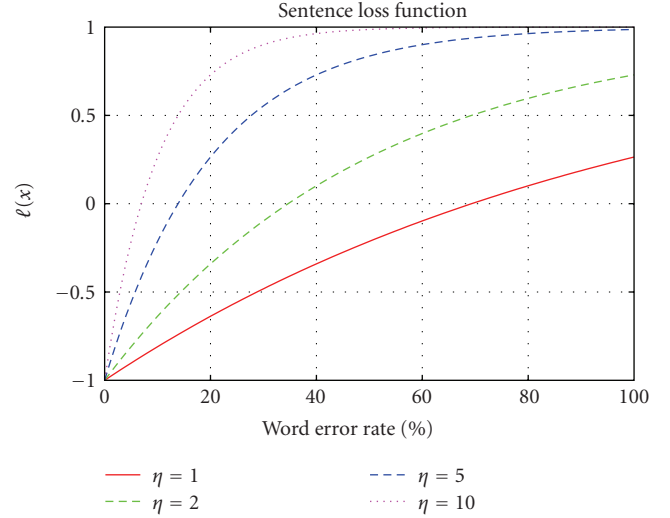


FIGURE 9: The sentence loss function (14) for  $\eta \in \{1, 2, 5, 10\}$ .

context. The following section discusses some cases for the distribution of  $y$  which are of relevance to the problem of speech recognition.

**6.1.3. Error Distributions in Sequential Decision Making.** In sequential decision-making problems, the knowledge about the correctness of decisions is delayed. Furthermore, it frequently lacks detailed information concerning the temporal location of errors. A common such case is knowing that we have made one or more errors in the time interval  $[1, T]$ . This form occurs in a number of settings. In the setting of individual sentence recognition, a sequence of decisions is made which corresponds to an inferred utterance. When this is incorrect, there is little information to indicate, where mistakes were made.

In such cases it is necessary to define a model (Even if no model is explicitly defined, there is always one implied.) for the probability of having made an erroneous decision at different points in times  $t$ , given that there has been at least one error in the interval  $[1, T]$ . Let us denote the probability of having made an error at time  $t \in [1, T]$  by  $\mathbf{P}(y_t \neq h_t \mid y_{1:T} \neq h_{1:T})$ . A trivial example of such a model is to assume that the error probability is uniformly distributed. This can be expressed via the flat prior

$$\mathbf{P}(y_t \neq h_t \mid y_{1:T} \neq h_{1:T}) \propto \frac{1}{T}. \quad (15)$$

Another useful model is to assume an exponential prior such that

$$\mathbf{P}(y_t \neq h_t \mid y_{1:T} \neq h_{1:T}) \propto \lambda^{t-T}, \quad \lambda \in [0, 1), \quad (16)$$

such that the expectation of an error near the end of the decision sequence is much higher. This is useful in tasks where it is expected that the decision error will be temporally close to the information that an error has been made. Ultimately, such models incorporate very little knowledge about the task, apart from this simple temporal structure.

In this case we focus on the application of speech recognition, which has some special characteristics that can be used to more accurately estimate possible locations of errors. For the case of labelled sentence examples, it is possible to have a procedure that can infer the location of an error in time. This is because correctly recognised words offer an indication of where possible errors lie. Assume some procedure that creates an indicator function  $I_t$  such that  $I_t = 1$  for instances in time, where an error could have been made. We can then estimate the probability of having an error at time  $t$  as follows:

$$\mathbf{P}(y_t \neq h_t \mid y_{1:T} \neq h_{1:T}) = \frac{\gamma^{I_t}}{\sum_{k=1}^T \gamma^{I_k}}, \quad (17)$$

where the parameter  $\gamma \in [1, \infty)$  expresses our confidence in the accuracy of  $I_t$ . A value of 1 will cause the probability of an error to be the same for all moments in time, irrespective of the value of  $I_t$ , while, when  $\gamma$  approaches infinity, we have absolute confidence in the inferred locations. Similar relations can be defined for an exponential prior, and they can be obtained through the convolution of (16) and (17).

In order to apply boosting to temporal data, where classification decisions are made at the end of each sequence, we use a set of weights  $\{\psi_t^{(i)}\}_i$  corresponding to the set of frames in an example sentence. At each boosting iteration  $j$  the weights are adjusted through the use of (17), resulting in the following recursive relation:

$$\psi_t^{(j+1)} = \frac{\psi_t^{(j)} \gamma^{I_t}}{\sum_{k=1}^T \psi_k^{(j)} \gamma^{I_k}}. \quad (18)$$

In this manner, the loss incurred by the whole sentence is distributed to its constituent frames, although the choice is rather ad hoc. A different approach was investigated by Zhang and Rudnicky [24], where the loss on the frames was related to the probability of the relevant word being uttered at time  $t$ , but their results do not indicate that this is a better choice compared to the simpler utterance-level training scheme that they also propose in that paper.

**6.2. Experiments with Expectation Boosting.** We experimented on the OGI Numbers 95 (N95) data set [28] (details about the setup and dataset are given in Section 4). The experiment was performed as follows: firstly, a set of HMMs  $\mu_0$ , composed of one model per phoneme, was trained using the available phonetic labels. This has the role of a starting point for the subsequent expert models. At each boosting iteration  $t$  we take the following steps: firstly, we sample with replacement from the distribution of training sentences given by the AdaBoost algorithm. We create a new expert  $\mu_t$ , initialised with the parameters of  $\mu_0$ . The expert is trained on the sentence data using EM with the Viterbi approximation in the expectation step to calculate the expectation. The frames of each sequence carry an importance weight  $\psi_t$ , computed via (18), which is factored into the training algorithm by incorporating it in the posterior probability of

the model  $h$  given the data  $x$  and time  $t$ , which, if we assume independence of  $x$  and  $t$ , can be written as

$$p(h \mid x, t) \propto p(h \mid x) p(h \mid t) p(h). \quad (19)$$

In this case,  $p(h \mid t)$  will correspond to our importance weight  $\psi_t$ .

After training, all sequences are decoded with the new expert. The weights of each sentence is increased according to (14), with  $\eta = 10$ . This value was chosen so that any sentence decodings with more than 50% error rate would be considered nearly completely erroneous (see Figure 9). For each erroneously decoded sentence we calculate the edit distance using a shortest path algorithm. All frames for which the inferred state belonged to one of the words that corresponded to a substitution, insertion, or deletion are then marked. The weights of marked frames are adjusted according to (17). The parameter  $\gamma$  corresponds to how smooth we want the temporal credit assignment to be.

In order to evaluate the combined models we use the multistream method described in (6), where the weight  $w_i$  of each stream  $i$  is  $w_i \triangleq \pi(h_i) = \beta_i / \sum_j \beta_j$ .

Experimental results comparing the performance of the above techniques to that of an HMM using segmentation information for training are shown in Figure 10(a) for the training data and Figure 10(b) for the test data. The figures include results for our previous results with boosting at the phoneme level. We have included results for values of  $\gamma \in \{1, 2, 4, 8, 16\}$ . Although we do not improve significantly upon our previous work with respect to the generalisation error, we found that on the training set, while boosting with presegmented phoneme examples had previously resulted in a reduction of the error to 3% after approximately 30 iterations (not shown), the sentence example training, combined with the error probability distribution over frames, converged to the same error after approximately 6 iterations. The situation was similar in the hold-out set, with the new approach converging to a good generalisation error at 10 iterations, while the previous approach required 16 iterations to reach the same performance. A drawback of the new approach, however, is the need to specify two new hyperparameters: (a) the shape of the cost function and (b) the shape of the expected error distribution. As mentioned previously, for (a) we are using (14) with  $\eta = 10$  and for (b) we have chosen a mix between a uniform distribution and an indicator function, with  $\gamma$  being a free parameter. Choosing  $\gamma$  is not trivial (i.e., it cannot be chosen in the training set), since large values can lead to overfitting, while values that are too small seem to provide no benefit. For the experiments described in Section 7 we will be holding out part of the training set in order to select an appropriate value for  $\gamma$ .

The main interesting feature of the utterance-level approach is that we are minimising the word error rate directly, which is the real objective. Secondly, we do not require segmentation information during training. Lastly, the temporal probability distribution, derived from the word errors and the state inference, provides us with a method to assign weights to parts of the decoded sequence. Its importance becomes obvious when we compare the performance



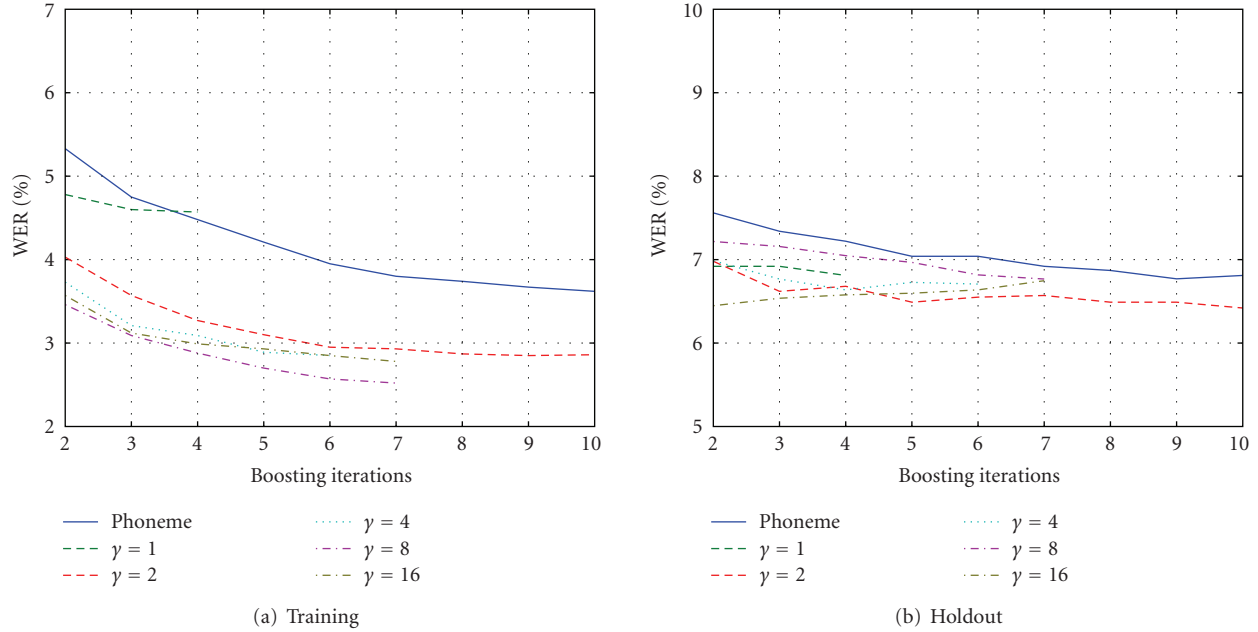


FIGURE 10: Word error rates for various values of  $\gamma$ , compared with the phoneme boosting approach, for the training and the holdout set.

of the method for various values of  $\gamma$ . When the distribution is flat (i.e., when  $\gamma = 1$ ), the performance of the model drops significantly. This supports the idea of using a probabilistic model for the errors over training sentences.

## 7. Generalisation Performance Comparison

In a real-world application one would have to use the training set for selecting hyperparameters to use in unknown data. To perform such a comparison between methods, the training data set was split in two parts, holding out 1/3rd of it for validation. For each training method, we selected the hyperparameters  $\theta$  having the best performance on the hold-out set. In our case, the hyperparameters are a tuple  $\theta = (k, n, \gamma)$ , where  $k \in \{10, 20, 30, 40, 50\}$  is the number of Gaussians in the Gaussian mixture model,  $n \in \{1, 2, 3, \dots, 16\}$  is the number of experts, and  $\gamma \in \{1, 2, 4, 8, 16\}$  is the temporal credit assignment coefficient for the expectation-boosting method. The number of states was fixed to 3, since our exploratory experiment, described in Section 5, indicated that it was the optimal value by a large margin. For each method, the hyperparameters  $\theta$  which provided the best performance in the hold-out set were used to train the model on the full training set. We then evaluated the resulting model on the independent test set. Full details on the data and methods setup are given in Section 4.

We compared the following approaches: firstly, a Gaussian mixture model (GMM), where the same observation distribution was used for all three states of the underlying phoneme hidden Markov model. This model was trained on the segmented data only; secondly, a standard hidden Markov model (HMM) with three states per phoneme, also trained on the segmented data only. We also considered the same models trained on complete utterances, using

*embedded training* after initialisation on the segmented data. The question we wanted to address was whether (and which) ensemble methods could improve upon these baseline results in an unbiased setting. We first considered ensemble methods trained using segmented data, using the phoneme-level bagging and boosting described in Section 5. This included both bagging and boosting of HMMs, as well as boosting of GMMs, for completeness. In all cases the experimental setup was identical, and the only difference between the boosting and the bagging algorithms was that bagging used a uniform distribution for each bootstrap sample of the data and uniform weights on the expert models. Finally, at the utterance level, we used expectation boosting, which is described in Section 6.

Table 1 summarises the results obtained, indicating the number of Gaussians per phoneme and the word error rate obtained for each model. If one considers only those models that were created strictly using the classification task, that is, without adapting word models, ensemble methods perform significantly better. Against the baseline *HMM embed* model, which is trained on full utterances, however, not all ensemble methods perform so well.

This can be seen clearly in Figures 11 and 12 which show the probability (Calculated via the bootstrap estimate, detailed in Section 4.2.) that one model will be better than the other by a given amount. In particular, the estimated probability that *Boost* is better than *HMM embed*, shown in Figure 12(a), is merely 51%, and the mean difference in performance is just 0.23% while, against the simple *HMM* the result, shown in Figure 11(a), is statistically significant with a confidence of 91%. Slightly better performance is offered by *E-Boost*, with significance with respect to the *HMM* and *HMM embed* models at 98% and 65%, respectively. Overall bagging works best, performing better

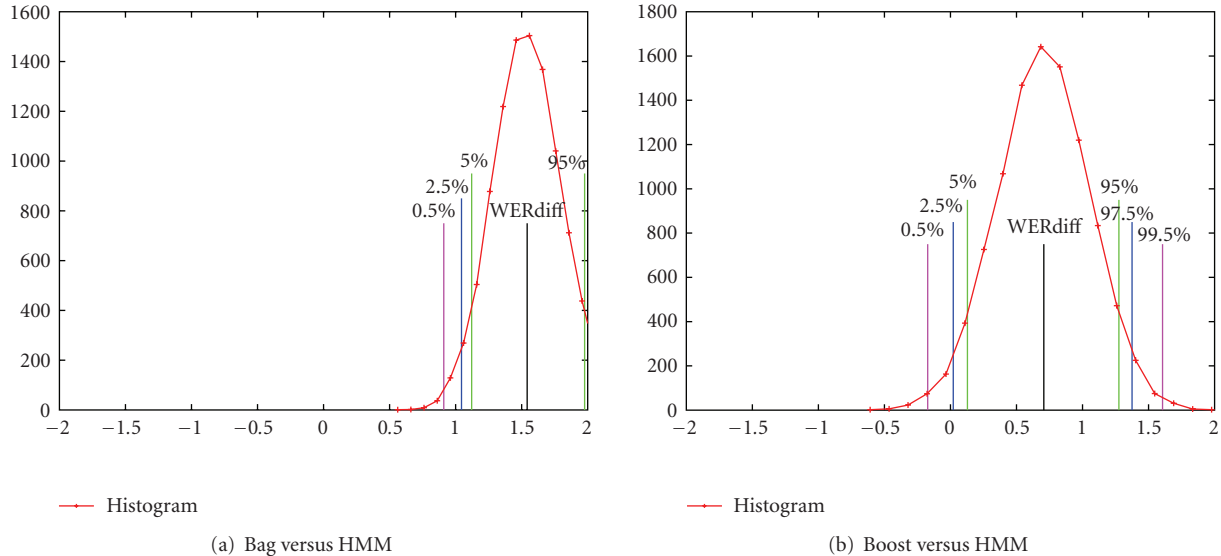


FIGURE 11: Significance levels of word error rate difference between the phoneme-level bagging and boosting and HMM. The histograms are created from 10,000 bootstrap samples of the test data, as described in Section 4.2. *WERdiff* displays the average difference in word error rate performance between the two methods. Positive values indicate that the first method has lower error than the second. The percentage values indicate the tail mass of the histogram to that point.

than other methods with a confidence of at least 99% in all cases, while approximately 97.5% of the probability mass lying above the 0.5% differential word error rate compared to the baseline model, as can be seen in Figure 12(a).

However, these results are not quite near the state of the art on this database. Other researchers (see, e.g., [36–40]) have achieved word error rates  $5.0 \pm 0.3\%$ , mainly through the use of different phonetic models. Accordingly, some further experiments were performed with Markov models using a more complex phonetic model (composed of 80 triphones, i.e., phonemes with contextual information). After performing the same model selection procedure as above, a single such model achieved word error rates of  $4.8 \pm 0.1\%$  (not shown in the table) which is in agreement with published state-of-the-art results. This suggests that using a more complex model could be better than using mixtures of simpler models. Corresponding results for ensembles of triphone models indicated that the boosting-based approaches could not increase generalisation performance, achieving a word error rate of 5.1%. However, the simpler bagging approach managed to reach a performance of 4.5%. However, the performance differences are not really significant in this case.

Nevertheless, it appears that, in all cases, phoneme bagging is the most robust approach. The reasons for this are not apparent, but it is tempting to conclude that the label noise combined with the variance-reducing properties of bagging is at least partially responsible. Although it should be kept in mind that the aforementioned triphone results are limited in significance due to the small difference in performance between methods, they nevertheless indicate that in certain situations ensemble methods and especially bagging may be of some use to the speech recognition community.

TABLE 1: Test set performance comparison of models selected on a validation set. The second column indicates the number of Gaussians per phoneme. For ensemble methods,  $n \times m$  denotes  $n$  models, each having  $m$  Gaussian components per state. *GMM* indicates a model consisting of a single Gaussian mixture for each phoneme. *HMM* indicates a model consisting of three Gaussian mixtures per phoneme. Thus, for HMMs, the total number of Gaussians is three times that of the GMMs with an equal number of components per state. *Boost* and *Bag* models indicate models trained using the standard boosting and bagging algorithm, respectively, on the phoneme classification task, while *E-boost* indicates the expectation boosting algorithm for word error rate minimisation. Finally *embed* indicates that embedded training was performed subsequently to initialisation of the model.

Model	Gaussians	Word error rate (%)
GMM	30	8.31
GMM embed	40	8.12
Boost GMM	$10 \times 30$	7.41
HMM	10	7.52
HMM embed	10	7.04
Boost HMM	$10 \times 10$	6.81
E-Boost HMM	$7 \times 10$ ( $\gamma = 8$ )	6.75
Bag HMM	$16 \times 20$	5.97

## 8. Discussion

We presented some techniques for the application of ensemble methods to HMMs. The ensemble training was performed for complete HMMs at either the phoneme or the utterance level, rather than at the frame level. Using boosting techniques at the utterance level was thought to lead

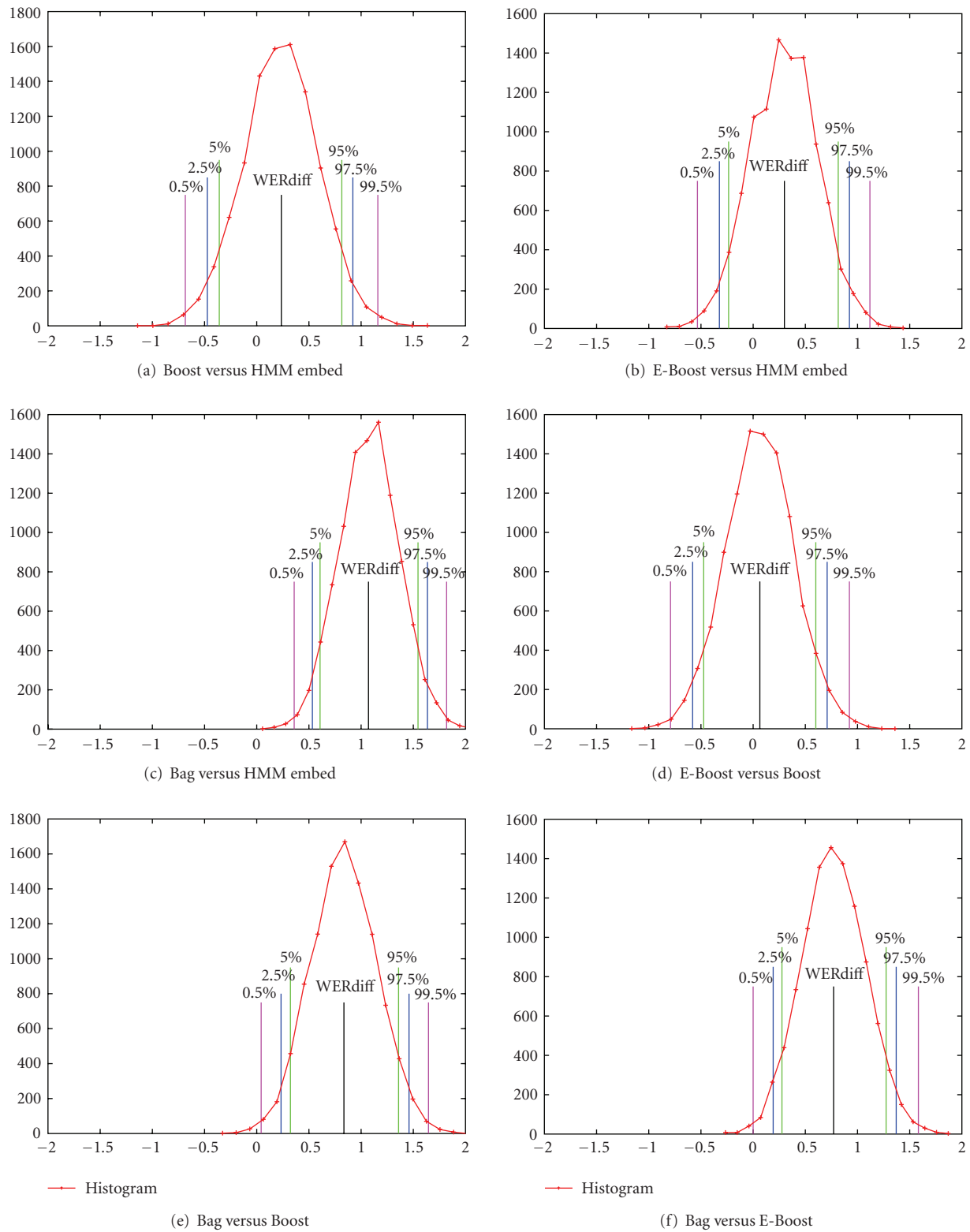


FIGURE 12: Significance levels of word error rate difference between the top four models. The histograms are created from 10,000 bootstrap samples of the test data, as described in Section 4.2. *WERdiff* displays the average difference in word error rate performance between the two methods. Positive values indicate that the first method has lower error than the second. The percentage values indicate the tail mass of the histogram to that point.

to a method for reducing the word error rate. Interestingly, this word error rate reduction scheme did not improve generalisation performance for boosting, while the simplest approach of all, bagging, performed the best.

There are a number of probable causes. The first one is that the amount of data is sufficiently large for ensemble techniques to have little impact on performance; that is, there is enough data to train sufficiently good base models. The second is that the state-locked multistream decoding techniques that were investigated for model recombination led to an increase in generalisation error as the inference performed is very approximate. The third is that the boosting approach used is simply inappropriate. The first case must not be true, since bagging does achieve considerable improvements over the other methods. There is some evidence for the second case, since the GMM ensembles are the only ones that should not be affected by the multistream approximations and, while a more substantial performance difference can be observed, it nevertheless is not much greater. The fact that bagging's phoneme mixture components are all trained on samples from the same distribution of data and that it outperforms boosting is also in agreement with this hypothesis. This leaves the possibility that the type of boosting training used is inappropriate, at least in conjunction with the decoding method used, open.

Future research in this direction might include the use of other approximations for decoding than constrained multistream methods. Such an approach was investigated by Meyer and Schramm [13], where the authors additionally consider the harder problem of large vocabulary speech recognition (for which even inferring the most probable sequence of states in a single model may be computationally prohibitive). It could thus be also possible to use the methods developed herein for large vocabulary problems by borrowing some of their techniques. The first technique, also used in [22], relies on finding an  $n$ -best list of possible utterances, assuming that there are no other possible utterances and then fully estimating the posterior probability of the  $n$  alternatives. The second technique, developed by Schramm and Aubert [41], combines multiple pronunciation models. In this case each model arising from boosting could be used in lieu of different pronunciation models. Another possible future direction is to consider different algorithms. Both AdaBoost.M1, which was employed here, and AdaBoost.M2, are using greedy optimisation for the mixture coefficients. Perhaps better optimisation procedures, such as those proposed by Mason et al. [42], may result in an additional advantage.

## Acknowledgments

This work was supported in part by the IST program of the European Community, under the PASCAL Network of Excellence, IST-2002-506778, and funded in part by the Swiss Federal Office for Education and Science (OFES) and the Swiss NSF through the NCCR on IM2, and the EU-FP7 project IM-CLeVeR.

## References

- [1] C. Dimitrakakis, *Ensembles for sequence learning*, Ph.D. thesis, École Polytechnique Fédérale de Lausanne, 2006.
- [2] C. Dimitrakakis and S. Bengio, "Boosting HMMs with an application to speech recognition," in *Proceedings of IEEE International Conference on Acoustics, Speech, and Signal Processing*, pp. 621–624, May 2004.
- [3] C. Dimitrakakis and S. Bengio, "Boosting word error rates," in *Proceedings of IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP '05)*, vol. 5, pp. 501–504, 2005.
- [4] A. Morris, A. Hagen, H. Glotin, and H. Bourlard, "Multi-stream adaptive evidence combination for noise robust ASR," *Speech Communication*, vol. 34, no. 1-2, pp. 25–40, 2001.
- [5] H. Misra and H. Bourlard, "Spectral entropy feature in full-combination multi-stream for robust ASR," in *Proceedings of the 9th European Conference on Speech Communication and Technology*, pp. 2633–2636, Lisbon, Portugal, 2005.
- [6] H. Misra, H. Bourlard, and V. Tyagi, "New entropy based combination rules in HMM/ANN multi-stream ASR," in *Proceedings of IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP '03)*, vol. 2, pp. 741–744, Hong Kong, 2003.
- [7] H. Ketabdaz, H. Bourlard, and S. Bengio, "Hierarchical multistream posterior based speech recognition system," 2005, IDIAP-RR 25, IDIAP.
- [8] L. Breiman, "Bagging predictors," *Machine Learning*, vol. 24, no. 2, pp. 123–140, 1996.
- [9] Y. Freund and R. E. Schapire, "A decision-theoretic generalization of on-line learning and an application to boosting," *Journal of Computer and System Sciences*, vol. 55, no. 1, pp. 119–139, 1997.
- [10] R. Meir and G. Rätsch, "An introduction to boosting and leveraging," in *Advanced Lectures on Machine Learning*, vol. 2600 of *Lecture Notes in Computer Science*, pp. 118–183, 2003.
- [11] R. E. Schapire and Y. Singer, "Improved boosting algorithms using confidence-rated predictions," *Machine Learning*, vol. 37, no. 3, pp. 297–336, 1999.
- [12] C. Breslin, *Generation and combination of complementary systems for automatic speech recognition*, Ph.D. thesis, Cambridge University Engineering Department and Darwin College, 2008.
- [13] C. Meyer and H. Schramm, "Boosting HMM acoustic models in large vocabulary speech recognition," *Speech Communication*, vol. 48, no. 5, pp. 532–548, 2006.
- [14] X. Yang, M.-h. Siu, H. Gish, and B. Mak, "Boosting with antimodels for automatic language identification," in *Proceedings of the 8th Annual Conference of the International Speech Communication Association (Inter-Speech '07)*, pp. 342–345, 2007.
- [15] R. Zhang and A. I. Rudnicky, "Apply  $n$ -best list re-ranking to acoustic model combinations of boosting training," in *Proceedings of the 8th International Conference on Spoken Language Processing (ICSLP '04)*, pp. 1949–1952, 2004.
- [16] C. Breslin and M. J. F. Gales, "Directed decision trees for generating complementary systems," *Speech Communication*, vol. 51, no. 3, pp. 284–295, 2009.
- [17] O. Siohan, B. Ramabhadran, and B. Kingsbury, "Constructing ensembles of asr systems using randomized decision trees," in *Proceedings of IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP '05)*, pp. 197–200, 2005.



- [18] H. Schwenk, "Using boosting to improve a hybrid HMM/neural network speech recognizer," in *Proceedings of IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP '99)*, vol. 2, pp. 1009–1012, 1999.
- [19] G. Zweig and M. Padmanabhan, "Boosting Gaussian mixtures in an LVCSR system," in *Proceedings of IEEE International Conference on Acoustics, Speech, and Signal Processing*, pp. 1527–1530, June 2000.
- [20] G. Cook and A. Robinson, "Boosting the performance of connectionist large vocabulary speech recognition," in *Proceedings of the International Conference on Spoken Language Processing (ICSLP '96)*, vol. 3, pp. 1305–1308, Philadelphia, Pa, USA, October 1996.
- [21] L. Bahl, P. Brown, P. de Souza, and R. Mercer, "A new algorithm for the estimation of hidden Markov model parameters," in *Proceedings of the IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP '88)*, pp. 493–496, 1988.
- [22] R. Zhang and A. I. Rudnicky, "Comparative study of boosting and non-boosting training for constructing ensembles of acoustic models," in *Proceedings of the 8th European Conference on Speech Communication and Technology (Eurospeech '03)*, pp. 1885–1888, 2003.
- [23] J. G. Fiscus, "Post-processing system to yield reduced word error rates: Recognizer Output Voting Error Reduction (ROVER)," in *Proceedings of IEEE Workshop on Automatic Speech Recognition and Understanding*, pp. 347–354, December 1997.
- [24] R. Zhang and A. I. Rudnicky, "A frame level boosting training scheme for acoustic modeling," in *Proceedings of the 8th International Conference on Spoken Language Processing (ICSLP '04)*, pp. 417–420, 2004.
- [25] M. H. Siu, X. Yang, and H. Gish, "Discriminatively trained GMMs for language classification using boosting methods," *IEEE Transactions on Audio, Speech and Language Processing*, vol. 17, no. 1, Article ID 4740154, pp. 187–197, 2009.
- [26] M. Gales and S. Young, "The application of hidden Markov models in speech recognition," *Foundations and Trends R in Signal Processing*, vol. 1, no. 3, pp. 195–304, 2007.
- [27] R. Zhang, *Making an Effective Use of Speech Data for Acoustic Modeling*, Ph.D. thesis, Carnegie Mellon University, 2007.
- [28] R. A. Cole, K. Roginski, and M. Fanty, "The OGI numbers database," Tech. Rep., Oregon Graduate Institute, 1995.
- [29] L. R. Rabiner and B.-H. Juang, *Fundamentals of Speech Recognition*, PTR Prentice-Hall, 1993.
- [30] J. Mariéthoz and S. Bengio, "A new speech recognition baseline system for Numbers 95 version 1.3 based on Torch," 2004, IDIAP-RR 04-16, IDIAP.
- [31] B. Efron and R. J. Tibshirani, *An Introduction to the Bootstrap*, vol. 57 of *Monographs on Statistics & Applied Probability*, Chapman & Hall, 1993.
- [32] M. Bisani and H. Ney, "Bootstrap estimates for confidence intervals in ASR performance evaluation," in *Proceedings of IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP '04)*, vol. 1, pp. 409–412, May 2004.
- [33] A. Van der Vaart and J. Wellner, *Weak Convergence and Empirical Processes: With Applications to Statistics*, Springer, Berlin, Germany, 1996.
- [34] C. M. Bishop, *Neural Networks for Pattern Recognition*, Clarendon Press, Oxford, UK, 1995.
- [35] G. Rätsch, T. Onoda, and K. R. Müller, "Soft margins for AdaBoost," *Machine Learning*, vol. 42, no. 3, pp. 287–320, 2001.
- [36] M. Athineos, H. Hermansky, and D. P. Ellis, "LP-TRAP: linear predictive temporal patterns," in *Proceedings of the 8th International Conference on Spoken Language Processing (ICSLP '04)*, pp. 949–952, 2004.
- [37] M. M. Doss, *Using auxiliary sources of knowledge for automatic speech recognition*, Ph.D. thesis, École Polytechnique Fédérale de Lausanne, Computer Science Department, Lausanne, Switzerland, 2005, Thesis no. 3263.
- [38] H. Hermansky and S. Sharma, "TRAPs—classifiers of temporal patterns," in *Proceedings of the 5th International Conference on Speech and Language Processing (ICSLP '98)*, pp. 1003–1006, 1998.
- [39] H. Ketabdar, J. Vepa, S. Bengio, and H. Bourlard, "Developing and enhancing posterior based speech recognition systems," in *Proceedings of the 9th European Conference on Speech Communication and Technology*, pp. 1461–1464, Lisbon, Portugal, September 2005.
- [40] G. Lathoud, M. Magimai-Doss, B. Mesot, and H. Bourlard, "Unsupervised spectral subtraction for noise-robust ASR," in *Proceedings of IEEE Automatic Speech Recognition and Understanding Workshop (ASRU '05)*, pp. 189–194, December 2005.
- [41] H. Schramm and X. L. Aubert, "Efficient integration of multiple pronunciations in a large vocabulary decoder," in *Proceedings of IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP '00)*, vol. 3, pp. 1659–1662, 2000.
- [42] L. Mason, P. L. Bartlett, and J. Baxter, "Improved generalization through explicit optimization of margins," *Machine Learning*, vol. 38, no. 3, pp. 243–255, 2000.